

Two-sided solution of ODEs via a posteriori error estimates

B.S. DOBRONEC

Computing Center, Academy of Sciences, Siberian Branch, Akademgorodok, Krasnoyarsk, U.S.S.R.

Received 5 May 1987

Revised 2 December 1987

Abstract: Two-sided methods for solving ODEs are presented in the paper. The methods are based on a posteriori error estimates. It is shown that these methods may be successfully applied to stiff ODEs. An illustrative numerical example is given.

Keywords: Two-sided methods, a posteriori error estimates, interval analysis, ordinary differential equations.

1. Introduction

Two-sided methods and interval analysis methods are two classes of numerical methods which construct error estimates. (About interval analysis methods see the overview in [1].) Usually two-sided methods do not use the notion of interval analysis, and in spite of their simplicity they have some disadvantages, for instance: the obtained error estimates are valid only for sufficiently small step sizes of the mesh and they do not consider the round-off error.

The aim of this work is to combine simplicity of two-sided methods and guaranteed precision of interval analysis methods. For this purpose the a posteriori error estimation is more useful, and well developed for linear differential equations (for literature see [2]).

The proposed method is as follows: first one solves a given problem and some additional problems by some numerical method and their numerical solutions are interpolated by splines. Then a two-sided solution is sought as a linear combination of these splines.

Interval analysis is applied here only for computing defects of the spline solutions and necessary constants. Thanks to this the time of computations may be considerably reduced and precision may be improved.

The advantages of this approach are that for its realization the known numerical methods for solving ODEs and spline approximations are used.

This approach is applied here in particular for 'stiff' ODEs.

2. Formulation of the problem

In this paper we consider ODEs of the form:

$$\frac{dy}{dx} = f(x, y), \quad x \in (0, l), \quad (1)$$

$$y(0) = y_0. \quad (2)$$

Here

$$y(x) = (y_1(x), y_2(x), \dots, y_n(x))^T,$$

$$y_0 = (y_{01}, y_{02}, \dots, y_{0n})^T, \quad y_{0i} \in [\underline{y}_{0i}, \bar{y}_{0i}] = y_{0i}^I,$$

$$f(x, y) = (f_1(x, y), f_2(x, y), \dots, f_n(x, y)),$$

$$f_i(x, y) \in f_i^I(x, y) = [\underline{f}_i(x, y), \bar{f}_i(x, y)],$$

and

$$f_i \in C^m([0, l] \times \mathbb{R}^n), \quad m \geq 3.$$

Suppose that a solution of (1), (2) exists on $[0, l]$, is unique and $y_i \in C^{m+1}[0, l]$, $i = 1, \dots, n$.

Let $N \geq 1$ be an integer and let the nodes of the mesh:

$$\omega_h = \{0 = x_0 < x_1 < x_2 < \dots < x_N = l\}$$

be a partition of $[0, l]$, $h = \max_{i=0, \dots, N-1} (x_{i+1} - x_i)$.

3. Construction of a two-sided solution

In order to construct a two-sided solution at first we approximately solve the problem (1), (2) for some concrete representatives of f , y_0 and use, for example, a Runger–Kutta method. As a result we get the approximate solution $y_i^h(x)$, $i = 1, \dots, n$ in the nodes ω_h . With the help of (1), (2) approximated values of the derivatives $\dot{y}_i^h(x)$ may be computed in the nodes ω_h . Using these values, we construct Hermite cubic splines S_i through the points $y_i^h(x)$, $x \in \omega_h$, $i = 1, \dots, n$.

Later we shall use the next defects

$$\phi_i(x, S) = f_i(x, s) - dS_i/dx, \quad i = 1, \dots, n,$$

$$S = (S_1, \dots, S_n).$$

Let us consider two additional ODE systems the solutions of which will be used for construction of the two-sided solution.

$$du/dx = Wu + w, \quad x \in (0, l), \quad u(0) = 0 \quad (3)$$

and

$$dv/dx = Wv, \quad x \in (0, l), \quad v(0) = z. \quad (4)$$

The matrix $W = \{W_{ij}\}$ consists of elements

$$W_{ii} = \frac{\partial f_i}{\partial y_i}(x, S), \quad i = 1, \dots, n, \quad W_{ij} = \left| \frac{\partial f_i}{\partial y_j}(x, S) \right|, \quad i \neq j, \quad i, j = 1, \dots, n,$$

vector w has components $w_i = 1$, vector z has components $\frac{1}{2}(\bar{y}_{0i} - y_{0i})$. Then we solve the problems (3), (4) by a Runge–Kutta method and build splines S_i^u , S_i^v which approximate the numerical solutions u_i^h , v_i^h , $i = 1, \dots, n$. We will seek a two-sided solution of the form

$$y_i^I = S_i + [-1, 1] S_i^v + \alpha^I S_i^u, \quad (5)$$

where $\alpha^I = [-\alpha, \alpha]$ is a certain interval constant. In order to find a constant α let us introduce the next functions ϕ_i^I , f_{yij}^I :

$$\phi_i^I(x, S) = [\underline{\phi}_i, \bar{\phi}_i] = f_i^I(x, S) - dS_i/dx,$$

$$\frac{\partial f}{\partial x}(x, \theta) \in f_{yij}^I(x, \theta^I) = [\underline{f}_{yij}, \bar{f}_{yij}], \quad \theta \in \theta^I.$$

Further, let $\delta_i > 0$ be a priori given components for which next inclusion is valid

$$y_i(x) \in S_i(x) + [-1, 1] S_i^v + [-\delta_i, \delta_i] S_i^u.$$

Remark 1. The constants δ_i may be found roughly. They may be estimated for example by a first order Moore method [4].

Then we set

$$\eta_i^I = S_i(x) + [-1, 1] S_i^v(x) + [-\delta_i, \delta_i] S_i^u(x)$$

and define

$$\tilde{f}_{yii}(x) = \tilde{f}_{yii}(x, \eta^I), \quad i = 1, \dots, n,$$

$$\tilde{f}_{yij}(x) = \max(|\underline{f}_{yij}(x, \eta^I)|, |\bar{f}_{yij}(x, \eta^I)|), \quad i \neq j, \quad i, j = 1, \dots, n,$$

$$\Phi_i(x) = \max(|\underline{\phi}_i(x, S)|, |\bar{\phi}_i(x, S)|) - dS_i^v/dx + \sum_{j=1}^n \tilde{f}_{yij}(x) S_j^v(x),$$

$$\Psi_i(x) = dS^u(x)/dx - \sum_{j=1}^n f_{yij}(x) S_j^u(x).$$

For a justification of the two-sided estimations we consider

$$\alpha = \max_{\substack{x \in [0, l] \\ i = 1, \dots, n}} (\Phi_i(x)/\Psi_i(x), 0) \quad (6)$$

where the functions Φ_i , Ψ_i are defined above. In order to find α we may use interval analysis provided that

$$\Psi_i(x) \neq 0, \quad x \in [0, l']. \quad (7)$$

In Section 4 it will be shown that there always exists a constant $l' > 0$, such that (7) takes place. If $l' < l$ then we may extend the two-sided solution from point $x = l'$ to $x = l$.

Remark 2. In the linear case the functions f_{yij} depend only on x . So for finding α it is not needed to know η_i^I and δ_i .

Now we prove inclusion (5). Consider the next equation for the errors $y_i - S_i$. Indeed, the S_i satisfy the following equations

$$dS_i(x)/dx = f_i(x, S) - \phi_i(x, S).$$

After subtraction from (1) we get

$$d(y_i - S_i)/dx = \sum_{j=1}^n \frac{\partial f_i}{\partial y_j}(x, \eta^i)(y_j - S_j) + \phi_i,$$

where the η^i are some functions in the intervals with bounds y, S . Consider the auxiliary problem

$$d\epsilon/dx = V\epsilon + \zeta, \quad x \in (0, l'), \quad \epsilon(0) = \epsilon_0, \quad (8)$$

where matrix V is such that

$$\begin{aligned} V_{ii}(x) &\geq \frac{\partial f_i}{\partial y_i}(x, \eta^i), \quad i = 1, \dots, n, \\ V_{ij}(x) &\geq \left| \frac{\partial f_i}{\partial y_j}(x, \eta^i) \right|, \quad i, j = 1, \dots, n, \quad i \neq j, \\ \epsilon_i^0 &\geq |y_{0i} - S_i(0)|, \quad \zeta_i(x) \geq |\phi_i(x, S)|, \quad x \in (0, l). \end{aligned} \quad (9)$$

Theorem 1 [5]. *Let ϵ be a solution of the ODEs (8) and let the condition (9) fulfilled. Then*

$$\epsilon_i(x) \geq |y_i(x) - S_i(x)|, \quad x \in [0, l].$$

So we put $V_{ij} = \tilde{f}_{y_{ij}}$. Evidently the conditions (9) are fulfilled for the matrix V . Consequently, from the construction of α we get that

$$\begin{aligned} (d(S^v + \alpha S^u)/dx - V(S^v + \alpha S^u))_i \\ = (dS^v/dx - VS^v + \alpha(dS^u/dx - VS^u))_i \geq |\phi_i|. \end{aligned}$$

Therefore

$$S_i^v + \alpha S_i^u \geq |y_i - S_i|$$

and hence

$$y_i \in S_i + [-1, 1] S_i^v + \alpha^1 S_i^u, \quad x \in [0, l],$$

i.e. we obtain the desired two-sided estimation.

4. Estimates of the width of the two-sided solution

Consider the errors $\epsilon_i(x) = y_i(x) - y_i^h(x)$, $x \in \omega_h$ and assume, that the next inequality is satisfied

$$|\epsilon_i(x)| = |y_i(x) - y_i^h(x)| \leq Ch^k, \quad \forall x \in \omega_h. \quad (10)$$

Let S_i^T be the Hermite cubic splines interpolating y_i on ω_h . Then [3]

$$\|d^v(y_i - S_i^T)/dx^v\|_{L_\infty[0, l]} \leq Ch^{4-v} \|y_i\|_{W_\infty^4[0, l]}. \quad (11)$$

Theorem 2. Let S_i be the Hermite cubic splines interpolating y_i^h on ω_h . Then

$$\|d^\nu(y_i - S_i)/dx^\nu\|_{L_\infty[0,l]} \leq C(h^{4-\nu} \|y_i\|_{W_\infty^4[0,l]} + h^k), \quad \nu = 0, 1. \quad (12)$$

Proof. Consider

$$\begin{aligned} \|d^\nu(y_i - S_i)/dx^\nu\|_{L_\infty[0,l]} &\leq \|d^\nu(y_i - S_i^T)/dx^\nu\|_{L_\infty[0,l]} \\ &\quad + \|d^\nu(S^T - S)/dx^\nu\|_{L_\infty[0,l]}, \quad \nu = 0, 1. \end{aligned}$$

Then by (10)

$$|dy_i(x)/dx - \dot{y}_i^h(x)| \leq Ch^k, \quad \forall x \in \omega_y. \quad (13)$$

Represent $S_i^T - S_i$ on the interval $[x_j, x_{j+1}]$ in the form

$$S_i^T - S_i = a_0 + a_1t + a_2t^2 + a_3t^3, \quad t = x - x_j.$$

Then by (10), (13)

$$|a_0|, \quad |a_1| \leq Ch^k, \quad |a_2| \leq Ch^{k-1}, \quad |a_3| \leq Ch^{k-2}.$$

Hence

$$\|d^\nu(S_i^T - S_i)/dx^\nu\|_{L_\infty[0,l]} \leq Ch^k, \quad \nu = 0, 1. \quad (14)$$

and use (11), (14) to obtain the statements of Theorem 2. \square

Now we estimate the width of the two-sided solution for

$$f_i^I \equiv f_i, \quad y_0^I \equiv y_0. \quad (15)$$

Note, that the width of the two-sided solution, depends on α , S_i^u , S_i^v . From (6) and (15) the next estimates follows:

$$\alpha \leq C \max_{i=1, \dots, n} |\phi_i(x, S)|.$$

By use of Theorem 2 we can estimate

$$\|\phi_i\|_{L_\infty[0,l]} \leq C(h^3 \|y_i\|_{W_\infty^4[0,l]} + h^4 \max_{i=1, \dots, n} \|y_i\|_{W_\infty^4[0,l]} + h^k).$$

Hence

$$\alpha \leq C \max_{i=1, \dots, n} (h^3 \|y_i\|_{W_\infty^4[0,l]} + h^k).$$

The width ρ of the two-sided solution can be bounded in the following way:

$$\rho(x) \leq C \max_{i=1, \dots, n} (h^3 \|y_i\|_{W_\infty^4[0,l]} + h^k) S_i^u(x).$$

Assume, that f_{yij}^I , f^I are inclusion monotonic [4]. Note, that if $\delta_i = \alpha$, then the two-sided solution will have a smaller width.

Now it is proved that there exists an interval $[0, l']$ such that (7) holds. Consider the functions

$$\tilde{y}_i(x) = dS_i^u(x)/dx - \sum_{j=1}^n W_{ij} S_j^4(x).$$

Then we have from Theorem 2 that

$$\tilde{\Psi}_i(x) \geq 1 - Ch^3, \quad \forall x \in (0, l), \quad i = 1, \dots, n.$$

There exists a positive constant C (see [4]) such that

$$|\tilde{f}_{yij} - W_{ij}| \leq C \sum_{j'=1}^n (S_{j'}^v + \delta_{j'} S_{j'}^u).$$

Hence

$$|\Psi_i - \tilde{\Psi}_i| \leq C \sum_{j=1}^n \left(\sum_{j'=1}^n (S_{j'}^v + \delta_{j'} S_{j'}^u) \right) S_j^u.$$

Because of $S_i^u(0) = 0$, $i = 1, \dots, n$ and the fact that S_i^u interpolates the solution of (3), there exists a positive constant C such that

$$|S_i^u(x)| \leq Cx, \quad \forall i = 1, \dots, n.$$

So, we have the following lemma.

Lemma 1. *There exist h, l' such that*

$$\Psi_i(x) > 0, \quad \forall x \in (0, l'), \quad i = 1, \dots, n.$$

5. Model problem and numerical example

Let us consider the next model problem

$$\begin{aligned} dy/dx &= bx + c, \quad x \in (0, 1), \quad y(0) = y_0, \\ y_0 &\in y_0^I, \quad b \in b^I, \quad c \in c^I, \quad y_0 = \frac{1}{2}(\underline{y}_0 + \bar{y}_0). \end{aligned} \quad (16)$$

Its exact solution is

$$y = ((y_0 b + c) \exp(bx) - c)/b.$$

Now we write the problems (3), (4) in the next form

$$du/dx = bu + 1, \quad x \in (0, 1), \quad u(0) = 0, \quad (17)$$

$$dv/dx = bv, \quad x \in (0, 1), \quad v(0) = \frac{1}{2}(\bar{y}_0 - \underline{y}_0) = \epsilon_0. \quad (18)$$

Let S, S^u, S^v be splines which interpolate the numerical solution of the problems (16)–(18). Then

$$\begin{aligned} \Phi(x) &= |b^I S(x) + c^I - dS(x)/dx| - dS^v(x)/dx + bS^v(x), \\ \Psi(x) &= -bS^u(x) + dS^u(x)/dx. \end{aligned}$$

Hence for sufficiently small h , we have from Lemma 1 the following estimate:

$$\Psi(x) \geq 1 - Ch^3 > 0.$$

So α exists and is defined by (6). Thus the two-sided solution has the next form:

$$y(x) \in S(x) + [-1, 1] S^v(x) + \alpha^I S^u(x).$$

The width ρ of the two-sided solutions is bounded in the following way:

$$\rho(x) \leq 2S^v(x) + 2\alpha S^u(x)$$

and

$$\rho(0) = \bar{y}_0 - \underline{y}_0.$$

Table 1

	1st component	2nd component	3rd component
$X = 0.001$			
upper bound	2.2411 E-8	1.0000	1.4710 E-3
exact solution	2.0611 E-9	0.9990	9.4907 E-4
lower bound	0.0	0.9808	4.2708 E-4
the width of the two-sided solution	2.2411 E-8	1.9571 E-3	1.0440 E-3
$X = 0.0002$			
upper bound	4.4120 E-17	0.9991	2.4687 E-3
exact solution	4.2484 E-18	0.9980	1.9461 E-3
lower bound	0.0	0.9970	1.4237 E-3
the width of the two-sided solution	0.4120 E-17	2.0600 E-3	1.0449 E-3
$X = 0.003$			
upper bound	1.7008 E-26	0.9981	3.4644 E-3
exact solution	8.7566 E-27	0.9970	2.9412 E-3
lower bound	0.0	0.9960	2.4183 E-3
the width of the two-sided solution	1.7008 E-26	2.1119 E-3	1.0460 E-3
$X = 0.004$			
upper bound	4.9175 E-34	0.9972	4.4581 E-3
exact solution	1.8049 E-35	0.9960	3.9342 E-3
lower bound	0.0	0.9905	3.4110 E-3
the width of the two-sided solution	4.9175 E-34	2.1578 E-3	1.0472 E-3
$X = 0.006$			
upper bound	3.0449 E-52	0.9953	6.4402 E-3
exact solution	7.6681 E-53	0.9940	5.9144 E-3
lower bound	0.0	0.9930	5.3879 E-3
the width of the two-sided solution	3.0449 E-52	2.2491 E-3	1.0522 E-3
$X = 0.008$			
upper bound	1.5837 E-69	0.9934	8.4144 E-3
exact solution	3.2576 E-70	0.9921	7.8871 E-3
lower bound	0.0	0.9910	7.3574 E-3
the width of the two-sided solution	1.5837 E-69	2.3352 E-3	1.0569 E-3

Since

$$\alpha \leq Ch^3, \quad S^u(x) \leq Cx, \quad S^v(x) \leq c_1 \epsilon_0 \exp(bx).$$

we have

$$\rho(x) \leq 2c_1 \epsilon_0 \exp(bx) + Ch^3x$$

and if $b \leq -Ch^3/(2\epsilon_0 c_1)$ then $\rho(x) \leq \rho(0)$.

Consider the stiff model problem arising in chemical kinetics [6]

$$\begin{aligned} dy_1/dx &= -20\,000 y_1, \\ dy_2/dx &= 20\,000 y_1 - y_2 + y_3, \\ dy_3/dx &= y_2 - y_3, \\ y_1(0) &= 1, \quad y_2(0) = y_3(0) = 0. \end{aligned} \tag{19}$$

The eigenvalues of the problem are $\lambda_1 = 0$, $\lambda_2 = -2$, $\lambda_3 = -20\,000$. In point $x = 0$ the solution of (19) has a boundary layer. For a construction of the two-sided solution at first we numerically solve the problem (19) by an implicit Runge–Kutta method.

The step $h_k = x_{k+1} - x_k$ of the mesh ω_h was variable $h_k \in [0.001, 0.1]$. Using this numerical solution we construct the special nonlinear splines defined in zone of boundary layer

$$R_i(x) = S_i + c_i^k \exp(\lambda_i^k x), \quad x \in [x_k, x_{k+1}], \quad k = 0, 1, \dots, N-1.$$

Here S_i are Hermite cubic splines, the constants C_i^k , λ_i^k are found by numerical solution y_i^h . Out of the boundary layer zone we use the Hermite cubic splines.

A two-sided solution was constructed in the next way. For every interval $[x_k, x_{k+1}]$, $k = 0, 1, \dots, N-1$ we solve the problems (3), (4) with a matrix W , where

$$W = \begin{bmatrix} -20\,000 & 0 & 0 \\ 20\,000 & -1 & 1 \\ 0 & -1 & 1 \end{bmatrix}.$$

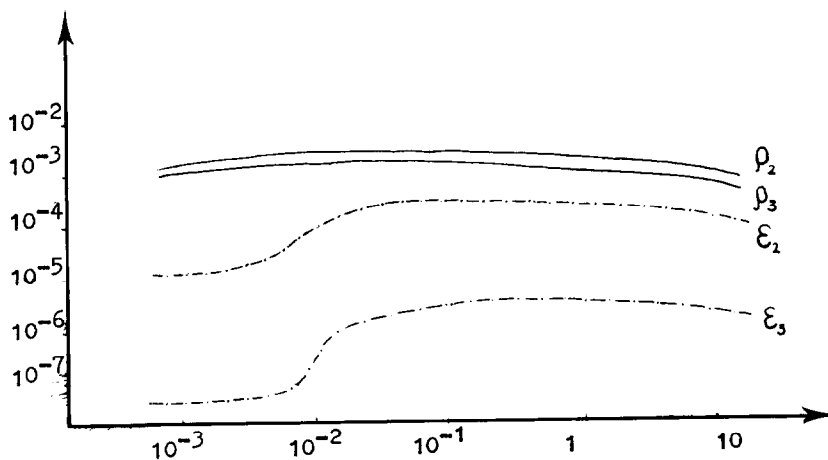


Fig. 1.

Because of linearity of the problem (19) the functions Φ_j , Ψ_j have the next form (see Remark 2)

$$\Phi_i = \sum_{j=1}^n W_{ij} (S_j + S_j^v) - d(S_i + S_i^v)/dx,$$

$$\Psi_j = dS_i^u/dx - \sum_{j=1}^n W_{ij} S_j^u.$$

By use of interval analysis we define by (6) the value of α for every interval $[x_k, x_{k+1}]$. Thus a two-sided solution has the form

$$y_i^1(x) = S_i(x) + [-1, 1] S_i^v(x) + [-\alpha_k, \alpha_k] S_i^u(x), \quad x \in [x_k, x_{k+1}].$$

Then the value $y_i^1(x_{k+1})$ is used as an initial value for solving the given problem on the interval $[x_{k+1}, x_{k+2}]$.

The computational results are shown in Table 1 and Fig. 1 where $\rho_i(x)$, $x \in \omega_h$, is the width of the two-sided solution, and $\epsilon_i(x) = |y_i(x) - y_i^h(x)|$, $x \in \omega_h$, is the error of the numerical solutions. All computations were performed on Es 1052 computer.

One can see that at the beginning width increases, and later stabilizes and decreases.

References

- [1] K. Nickel, Using interval methods for the numerical solution of ODEs, *ZAMM* **66** (1986) 513–523.
- [2] B.S. Dobronec, Dvustoronnie metody reshenija nekotorych uravnenij matfiziki, Dis. Kand. Fyz., Mat. Nauk, Novosibirsk, 1985.
- [3] R.S. Varga, Functional analysis and approximation theory in numerical analysis, SIAM, 1971.
- [4] R.E. Moore, *Interval Analysis* (Prentice-Hall, Englewood Cliffs, NY, 1966).
- [5] S.M. Lozinskij, Ocenka pogrešnosti približennogo reshenija sistemy differencial'nych uravnenij, *DAN SSSR* **92** (1953) 225–228.
- [6] V.I. Bykov and B.S. Dobronec, Dvustoronnie metody reshenija uravnenij himičeskoj kinetiki, V. sb.: Čislennye metody mehaniki splošnoj sredy, Novosibirsk: VC SO AN SSSR, 1985, t. 16, No. 4, pp. 13–22.